April 12, 2023

# RFP

# Technical Requirements Document

## for

# NERSC-10 System

Version 1.0

# TABLE OF CONTENTS

# REQUIREMENTS DEFINITIONS

Technical requirements have priority designations, which are defined as follows:

**(a)      Mandatory Requirements designated as (MR)**

Mandatory Requirements (designated MR) are performance features that are essential to NERSC-10 requirements, and an Offeror must satisfactorily propose all Mandatory Requirements in order to have its proposal considered responsive.

**(b)      Mandatory Option Requirements designated as (MO)**

Mandatory Option Requirements (designated MO) are features, components, performance characteristics, or upgrades whose availability as options are mandatory, and an Offeror must satisfactorily propose all Mandatory Option Requirements in order to have its proposal considered responsive. The Laboratory may or may not elect to include such options in the resulting subcontract(s). Therefore, each MO shall appear as a separately identifiable item in Offeror's proposal.

**(c)      Target Requirements designated as (TR-1, TR-2 or TR-3)**

Target Requirements (designated TR-1, TR-2, or TR-3) are features, components, performance characteristics, or other properties that are important to the Laboratory, but that will not result in a non-responsive determination if omitted from a proposal. Target Requirements are prioritized by dash number. TR-1 is most desirable to the Laboratory and forms the baseline system, while TR-2 is more desirable and adds additional capabilities or increases productivity. TR-3s are stretch goals. Target Requirement responses will be considered as part of the proposal evaluation process.

**(d)      Technical Option Requirements designated as (TO-1, TO-2 or TO-3)**

Technical Option Requirements (designated TO-1, TO-2, or TO-3) are features, components, performance characteristics, or upgrades that are important to the Laboratories, but that will not result in a non-responsive determination if omitted from a proposal. Technical Options add value to a proposal. Technical Options are prioritized by dash number. TO-1 is most desirable to the Laboratory, while TO-2 is more desirable than TO-3. Technical Option responses will be considered as part of the proposal evaluation process; however, the Laboratory may or may not elect to include Technical Options in the resulting subcontract(s). Each proposed TO should appear as a separately identifiable item in an Offeror's proposal response.


*Note: There are no mandatory requirements or mandatory options in the NERSC-10 technical requirements document.*

# 1.0 INTRODUCTION

The Regents of the University of California (the "University"), which operates the National Energy Research Scientific Computing ("NERSC") Center residing within Lawrence Berkeley National Laboratory ("LBNL"), is releasing a Request for Proposal (RFP) for the next generation high performance computing (HPC) system, NERSC-10 to be delivered in the 2026 time frame.

The successful NERSC-10 Offeror will be responsible for delivering, installing, supporting and maintaining the NERSC-10 System.

Each response/proposed solution within this document shall clearly describe the role of any lower-tier subcontractor(s) and the technology or technologies, both hardware and software, and value added that the lower-tier subcontractor(s) provide(s), where appropriate.

The scope of work and technical specifications for any subcontracts resulting from this RFP will be negotiated based on this Technical Requirements Document and the successful Offeror's responses/proposed solutions.

NERSC-10 has maximum funding limits over its system life, to include all design and development, site preparation, maintenance, support and analysts. Total ownership costs will be considered in system selection. The Offeror must respond with a configuration and pricing.

Application performance and workflow efficiency are essential to these procurements. Success will be defined as meeting NERSC-10 mission needs. The advanced workflows aspects of the NERSC-10 system will be pursued both by fielding first of a kind workflow enabling technologies as part of the system and by selecting and participating in strategic Non-Recurring Engineering (NRE) projects with the Offeror and applicable technology providers. A compelling set of NRE projects will be crucial for the success of NERSC-10, by enabling the deployment of first-of-a-kind technologies in such a way as to maximize their utility.

Supporting information can be found on the NERSC-10 website (https://www.nersc.gov/systems/nersc-10/).

Additional information on proposal preparation will be provided in the ***Proposal Submittal Requirements Section of the RFP***.

## 1.1 NERSC-10 Mission Need

The DOE Office of Science (SC) is the largest supporter of basic and applied research programs in the areas of efficient energy use, reliable energy sources, improved environmental quality, and fundamental understanding of matter and energy. One of the principal thrusts within SC is the direct support of the development, construction, and operation of unique, open-access High Performance Computing (HPC) scientific user facilities. These HPC facilities are critical to supporting the research programs that help accomplish the DOE's mission.

The [National Energy Research Scientific Computing](#) (NERSC) Facility at Lawrence Berkeley National Laboratory, funded by the DOE SC's [Advanced Scientific Computing Research](#) (ASCR) Office, is the mission HPC facility for the Office of Science, uniquely supporting the needs of science across the entire office. ASCR's mission is to discover, develop, and deploy computational and networking capabilities to analyze, model, simulate, and predict complex phenomena important to the DOE. The ASCR program has a long history supporting cutting edge research in applied math, computational and computer science and the deployment of advanced HPC and networking facilities.
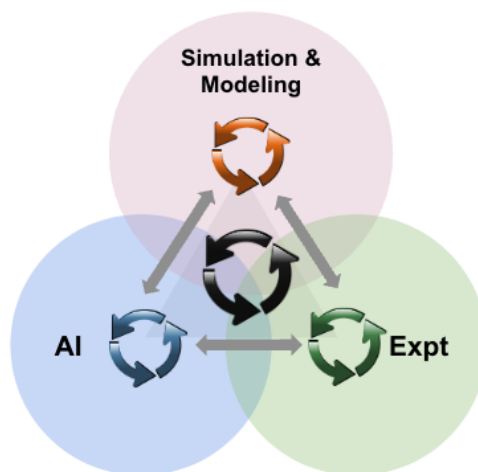
For [almost 50 years NERSC has supported SC](#) (and its predecessor agencies) funded researchers using cutting-edge supercomputers to develop materials and devices for clean energy generation and storage, study the environmental impact of a changing global ecosystem, investigate the fundamental properties and interactions of matter, and explore other science areas within the DOE science mission. Over that time NERSC has helped guide the SC computational science community through many disruptive changes and evolving national scientific priorities.

The SC community is now witnessing another transition with new science use cases requiring more flexible and programmable HPC systems, which presents new opportunities for mapping computational tasks and workflow components to technologies that might be more efficient and targeted for the workload. For SC to fulfill its mission to maintain and extend U.S. leadership in scientific discovery, the NERSC-10 system must leverage available new technologies and support the emerging needs in AI and experimental/observational science by accelerating end-to-end DOE SC workflows and enabling new modes of scientific discovery through the integration of experiment, data analysis, and simulation.

A NERSC-10 system upgrade in the 2026 timeframe is essential to meet Office of Science goals and national initiatives, and to avoid creating a gap between SC programmatic needs and HPC capabilities. The gap goes beyond the well-documented and nearly insatiable demand for more computational cycles from traditional simulation and modeling workloads. The NERSC-10 system must also address the growing need for more dynamic and programmable systems to accommodate increasingly complex workflows that may require running many interdependent simulations and/or analysis tasks, moving vast amounts of complex data among storage hierarchies both within NERSC and externally, for providing interactive real-time feedback to experiments, and for integrating simulations with AI.

## 1.2 HPC Workflows

The NERSC workload is increasingly diverse, with a growing demand for complex high performance workflow capabilities from our user community. NERSC supports and will continue to support HPC science campaigns through

- High-performance simulation and modeling workflows (e.g., large-scale multi-physics applications)
- High performance AI workflows (e.g., training, inference, hyperparameter optimization)
- Experiment to high-performance data analytics workflows (e.g., data analysis)

The new HPC workflow couples the computational tasks and data flow within, between or across all three modes to unlock opportunities for new scientific discovery. Enabling these complex and novel workflows requires increases in high-performance computing, composable resources for targeted workload performance, modular workflow execution through orchestration frameworks and APIs, spanning resources into commercial cloud, other computational facilities and the edge, tight integration of system components to enable 'seamless' execution, leveraging data stored/collected at other facilities, and more.

The NERSC-10 system must enable numerous simultaneous workflows with different resource and time requirements. NERSC is interested in technologies that can help support this vision. The **NERSC-10 Workflows Whitepaper** provides workflow scenarios that we expect to support with the NERSC-10 system.

## 1.3 Schedule

The following is the tentative schedule for the NERSC-10 system.

**Table 1. NERSC-10 high-level schedule.**

|  | NERSC-10 |
|---|---|
| RFP Released | Q4 CY 2023 |
| Subcontracts (NRE/Build) Awarded | Q4 CY 2024 |
| Test or Early Access System (Pilot or Phase I) | CY2025 |
| On-site System Delivery | 2H CY2026 |
| Production | CY2027 |

RFP Technical Requirements Document for NERSC-10 System, Version 1.0, April 12, 2023

# 2.0 HIGH-LEVEL SYSTEM REQUIREMENTS

This section describes the high-level technical requirements for the NERSC-10 system proposals. In addition, the RFP provides system and node tables for the Offeror to complete and to submit separately as well as to include here.

In addition to the TRs identified in this document, the Offeror may choose to propose any additional features (i.e. Offeror-proposed features) consistent with the objectives of the NERSC-10 procurement and the Offeror's roadmap, which the Offeror believes will be of value to the University.

## 2.1 System Description

2.1.1 The NERSC-10 system shall be sited at National Energy Research Scientific Computing Center data center in Building 59 on the Lawrence Berkeley National Laboratory campus in Berkeley, California. The Offeror should provide details of the physical footprint of the system and all of the supporting components to be sited at the NERSC data center to meet the facility requirements in Section 9. [TR-1]

2.1.2 The Offeror should provide a detailed full system architectural description of the NERSC-10 system that will **deliver at least a 10x Workflow-SSI performance improvement over Perlmutter**, including diagrams and text describing the following details as they pertain to the Offeror's proposed system architecture(s) plus any unique features in the design. Include quantities and define any minimum scalable unit sizing to maintain optimal performance and productivity across the system. [TR-1]

- Component architecture – details of all processor(s), memory technologies, storage technologies, network interconnect(s) and any other applicable components.

- Compute node architecture(s) – details of how components are combined into the node architecture(s). Details shall include bandwidth and latency specifications (or projections) between components. Details should be provided for each compute node type in the system. NERSC-10 should contain both CPU-only and GPU nodes, the balance of CPU-only to GPU nodes should be chosen to optimize the Workflow-SSI described in Section 3.0.

- Board and/or blade architecture(s) – details of how the node architecture(s) is integrated at the board and/or blade level. Details should include all inter-node and inter-board/blade communication paths and any additional board/blade level components.

- Rack and/or cabinet architecture(s) – details of how board and/or blades are organized and integrated into racks and/or cabinets. Details should include all inter rack/cabinet communication paths and any additional rack/cabinet level components.

- Interconnect - details of the system's high speed network topology and connectivity across all system components (compute nodes, workflow environment nodes, Platform Storage, QoS Storage, management system).

- Storage Systems – details of how the Platform Storage and QoS Storage are integrated with the system, including an architectural diagram and gateway nodes if applicable.

- System architecture – details of how rack or cabinets are combined to produce system architecture, including the high-speed interconnects and network topologies (if multiple) and storage systems.

- Management node(s) - details of hardware to support management and services to operate the NERSC-10 system. Management node types can include, but are not limited to, master nodes for orchestration of system services, worker nodes for the deployment of services, Slurm resource manager, and storage nodes for system management and administration. Multiple node types may be needed to optimize for different uses described in Section 5.0.

- Workflow environment node(s) (WEN) - details of hardware to support user access and user-driven workflow activities. A pool of WENs will be needed to address the different requirements described in Section 4.0. NERSC-10 may require multiple WEN types optimized for different use cases.

2.1.2   The Offeror should provide an alternative processor vendor and technology. The response shall concisely describe a basic architectural description, including hardware and software.  The response shall include a description of the PCI-e technology generation, BabelStream and py-GEMM microbenchmark results and power estimates (See Section 3 for benchmark descriptions). The Offeror should describe any other technologies that are part of their NERSC-10 offering that are high risk and propose mitigation strategies, including alternative technologies. [TR-1]

## 2.2 Software Description

2.2.1   The Offeror should provide a detailed description of the proposed system software and user programming environment, including a high-level software architecture diagram, the provenance of the software component (for example open source or proprietary), support mechanism and licensing, if applicable (for the lifetime of the system including updates). [TR-1]

2.2.2   The Offeror should describe the high-level roadmap for: [TR-1]

- System software and tools provided for management and operation of the NERSC-10 system.
- Provided user programming environment, including the ability to utilize new hardware features.

## 2.3 Non-Recurring Engineering (NRE)

The University expects to award Non-Recurring Engineering (NRE) subcontracts, separate from the system build subcontract. It is expected that NERSC personnel will collaborate in NRE subcontracts. It is anticipated that the NRE subcontracts could be approximately 10%-20% of the NERSC-10 system budgets. The Offeror is encouraged to provide proposals for areas of collaboration they feel provide substantial value to the NERSC-10 system (e.g., TR-2 or TR-3 requirements). The goals of the NRE efforts include, but are not limited to:

- Optimizing the usage of new hardware and software
- Enhancing NERSC-10 for Workflows
  - Increasing performance and portability
  - Enabling composability, automation & seamlessness, debugging & profiling, and monitoring
- Enhancing the NERSC-10 system
  - Increasing resilience, reliability and security
  - Enabling seamless integration into data center
  - Optimizing power and energy usage

Proposed collaboration areas should focus on topics that provide added value beyond planned roadmap activities. Proposals should not focus on one-off point solutions or gaps created by their proposed design that should be otherwise provided as part of a vertically integrated solution.

## 2.4 Upgrades, Expansions and Additions

The University expects to have future requirements for system upgrades and/or additional quantities of components based on the configurations proposed in response to this solicitation. The Offeror should propose separately priced options using whatever is the natural unit for the proposed architecture design as determined by the Offeror. For example, for system size, the unit may be the number of racks or some other unit appropriate for incrementally increasing the system. The Offeror should identify any thresholds requiring increased component infrastructure (e.g., extra spine switches), any technical challenges foreseen with respect to scaling and any other production issues. Proposals should be as detailed as possible.

2.4.1 The Offeror should propose and separately price upgrades, expansions or procurement of additional system configurations by the following fractions of the system as measured by the Workflow Sustained System Improvement (Workflow-SSI) metric. [TO-1]

- 25%
- 50%
- 100%
- 200%

2.4.2 The Offeror should propose upgrades, expansions or procurement of additional platform storage system capacity and QoS storage system capacity in increments of 10% for the scalable units described in Section 7. [TO-1]

2.4.3 The Offeror should propose non-volatile storage options (NVMe and/or SSD) internal to all node types, with at least triple the node memory capacity. Non-volatile storage performance characteristics and system software requirements to access and manage the storage should be described. [TO-1]

2.4.4 The Offeror should propose a novel AI acceleration partition that will both accelerate AI workloads and integrate within an HPC ecosystem to enable workflow capabilities. The partition should be delivered with the main NERSC-10 system or later and be made up of the scalable unit for the proposed design as determined by the Offeror. The description shall include: [TO-3]

- an overall architectural diagram that shows all hardware, interconnect(s), compilation infrastructure, and I/O subsystems, if applicable.
- an overview of software architecture, including libraries and SDKs, support for frameworks (e.g., TensorFlow, PyTorch, etc.), usability and programmability. The description should include terms of software licensing and any support if applicable.
- available results or projections on MLPerf benchmarks, in particular the "Training" and "Training: HPC" and "Inference Datacenter" benchmark suites. Results and projections provided should specify the version of the benchmark used, whether the result was officially submitted and any required modifications to the benchmark rules required to obtain the reported results or projections.
- performance results (actual, predicted or extrapolated) for the proposed system for one or more of the workflow component benchmarks listed in Table 3.1.

## 2.5 Early Access Systems

To allow for early and/or accelerated development of applications or development of functionality required as a part of the statement of work, the Offeror should propose options for Early Access Systems (EAS). The early access systems should contain similar functionality to the final system, including storage systems, management and workflow environment nodes, but scaled down to the appropriate configuration.

2.5.1 The Offeror should propose a **NERSC-10 Phase 1 Early Access System** that could be delivered in 2025 for production use. The primary purpose is to expose applications to the same programming environment as will be found on the final system. It is acceptable for the early access system to not use the final processor, node, or high-speed interconnect architectures. However, the programming and runtime environment must be sufficiently similar that a port to the final system is straightforward. The Offeror should propose an option that is 10% of the performance of the final NERSC-10 system (equivalent performance to the Perlmutter system). [TO-1]

2.5.2   The Offeror should propose a small **"Pilot" Early Access System** that could be delivered in 2025 to aid in the integration and development of the NERSC-10 system, form the basis for collaborative engineering efforts, or reduce risk and/or accelerate development.  [TO-1]

2.5.3   The Offeror should propose other **Early Access Systems** (hardware and software) to aid in the integration and development of the NERSC-10 system, form the basis for collaborative engineering efforts, or reduce risk and/or accelerate development. Of particular interest are resources that are in support of any topics proposed for NRE and workflow readiness. The systems could be on site at NERSC and/or accessed remotely, and could include access to early software, simulators and/or emulators. [TO-1]

## 2.6 Test and Development Systems (TDS)

A test and development system (TDS) shall contain all the functionality of production NERSC-10 systems, including storage systems, all accelerator types, but scaled down to the appropriate configuration. Multiple TDSs may be awarded to aid with the lifetime system management of any production system. It is desirable for NERSC-10 production systems and TDSs to be able to dynamically attach and detach from the same resources to allow scale testing on the test system by temporarily moving these resources from the production system to the test system.  The Offeror should propose Test and Development Systems for any production system delivered in support of the NERSC-10 system (e.g., the Phase 1 System described in 2.5.1 and the main NERSC-10 system). The TDSs should be delivered before the production resource they are designed to support.

2.6.1   The Offeror should propose a Pre-production TDS, which should contain at least 32 compute nodes. [TO-1]

2.6.2   The Offeror should propose a System Development TDS, which should contain at least 16 compute nodes. [TO-1]

# 3.0 BENCHMARKS

Assuring that real workflows perform well on the NERSC-10 system is key to the success of the system. The workflow component benchmarks listed in Table 3.1 will be used to evaluate workflow performance as part of both the RFP response and system acceptance. These benchmarks demonstrate the most computationally intensive components of the workflows they represent, but without the data-flow and control-flow complexities of an integrated science workflow. The workflow benchmarks are supplemented by a collection of micro-benchmarks listed in Table-3.2.

Final benchmark acceptance performance targets will be negotiated after a final system configuration is defined. All performance tests must continue to meet acceptance criteria throughout the lifetime of the system.

The benchmarks and supplemental materials can be found on the NERSC-10 benchmarks website: https://www.nersc.gov/systems/nersc-10/benchmarks. All benchmark results shall be reported in the accompanying "NERSC-10 Benchmark Results" worksheet.  All benchmark results shall conform to the "Workflow Component Benchmark and Micro-benchmark Instructions and Run Rules" document.  Benchmark results can be submitted for three categories of workflow optimization (baseline, ported and optimized), which are defined in the Run Rules.

**Table 3.1. Workflow Component Benchmarks**

| Workflow Name | Description | Application Components |
|---|---|---|
| Lattice QCD | Lattice Quantum Chromodynamics (QCD) | MILC generation MILC analysis |
| Optical Properties of Materials | *Ab initio* Electronic Structure | BerkeleyGW Epsilon BerkeleyGW Sigma |
| Materials by Design | Molecular Dynamics | LAMMPS |
| Climate Simulation & Analysis | Deep Learning Training | DeepCAM |
| Metagenome Analysis | Genomic Data Analysis | HMMsearch |
| CMB-S4 | Cosmology Data Analysis | TOAST-3 |

3.0.1    The Offeror should provide baseline performance results for the proposed system and platform storage system for all of the workflow component benchmarks listed in Table-3.1. [TR-1]

3.0.2    The Offeror should provide ported performance results for the proposed system and platform storage system for any of the workflow component benchmarks.  [TR-2]

3.0.3    The Offeror should provide optimized performance results for the proposed system and platform storage system for any of the workflow component benchmarks. [TR-2]

3.0.4    The Offeror should state a minimum workflow-SSI for the NERSC-10 system, to be measured using baseline versions of the workflow component benchmarks. If baseline results cannot be obtained, ported results may be provided in their place. [TR-1]

3.0.5    The Offeror should state a minimum workflow-SSI for the NERSC-10 system, to be measured using any combination of baseline, ported or optimized versions of the workflow component benchmarks. [TR-2]

3.0.6 The Offeror should provide baseline performance results for the proposed compute system and QoS Storage System for the DeepCAM workflow component benchmark. [TR-2]

3.0.7 The Offeror should provide performance results for the proposed system for micro-benchmarks listed in Table 3.2. Some micro-benchmarks can be run on multiple subdomains of the system; results should be provided for each configuration listed in the table. [TR-1]

3.0.8 The Offeror should provide licenses for the delivered system for all software required to achieve benchmark performance, including, but not limited to compilers and libraries. [TR-1]

**Table 3.2. Micro-benchmarks**

| Name | Description | Test Configuration(s) |
|------|-------------|----------------------|
| BabelStream | Memory bandwidth | Every functional combination of processor type and memory domain |
| py-DGEMM | Floating-point performance | Every processor type |
| Ziatest | Job startup | Full-system |
| OSU Micro-Benchmarks (OMB) | MPI performance | Every functional combination of processor type and memory domain. It is not necessary to test communication between heterogeneous processor/ memory combinations. |
| iperf | External networking performance | HSN ↔ PSS<br>HSN ↔ QSS |
| IOR | Storage bandwidth | PSS - full<br>QSS - minimum scalable unit<br>QSS - multiple independent scalable units<br>QSS - full |
| MDTest | Storage metadata | PSS - full<br>QSS - minimum scalable unit<br>QSS- multiple independent scalable units<br>QSS - full |

# 4.0 WORKFLOW ENVIRONMENT

## 4.1 Scalable and Reliable Workflow Services

4.1.1    The system should support running jobs up to the full scale of the compute node resources. The Offeror should describe factors (such as executable size) that may affect application launch time. [TR-1]

4.1.2    SchedMD's Slurm resource job management scheduler will be the primary scheduler and policy engine of the system. The University will directly procure the necessary software licenses and ongoing maintenance support from SchedMD. The Offeror will work with the University and/or SchedMD to resolve operational problems with Slurm that may be caused by the Offeror's products. The Offeror will provide the necessary integration interfaces to support scalable job launch, including node placement, topology-aware scheduling, rank reordering, power-aware scheduling, and node configuration and re-provisioning of nodes if supported by the hardware. The system design should not limit Slurm's ability to support thousands of concurrent users and more than 20,000 concurrent batch jobs. [TR-1]

4.1.3    The system should support a container orchestration platform such as Kubernetes or similar to provide staff-managed and user self-supported services on the workflow environment nodes. It should be capable of operating with Slurm to provide a unified workflow environment in which users can securely and performantly launch job tasks on the compute resources from the workflow environment nodes or services running on them. The Offeror should describe any specialized hardware [or software] that may be required to support or enhance this unified workflow environment capability. [TR-2]

4.1.4    The system workflow environment nodes shall support interactive user access modes, including: [TR-1]

   ● command-line interface (CLI) through ssh and web-based user access modes for login, code compilation (cross-compilation is not desirable), application development, container builds, job lifecycle management, small-scale data analysis, and data transfer.
   ● long-lived user services and frameworks (e.g., JupyterHub, databases, API services, and message brokers), that are staff-managed or user self-supported.

The Offeror should describe mechanisms to enable these access models and how they are managed.

4.1.5    The system should provide correct numerical results and minimize runtime variability. The Offeror should describe strategies for minimizing runtime variability in production.  [TR-1]

## 4.2 Programming Environment and Software Tools

4.2.1   The system should support building and executing C17 code, C++20 code and Fortran 2018 code including code utilizing OpenMP directives 5.2 or latest. The Offeror should describe all supported compilers, including any enhancements or limitations that can be expected in meeting full support of the standards and other native language features for expressing parallelism including, but not limited to, support for C++ parallel STL, Fortran *do concurrent* and coarrays. [TR-1]

4.2.2   The system should support LLVM backends for each processing element (e.g., CPU, GPU, specialized accelerator) that can be utilized with both vendor provided frontends and the open clang/flang projects. [TR-2]

4.2.3   The Offeror should describe the capability of the system to compile and run applications using Kokkos, SYCL, and/or OpenACC 3.x. If present, describe any performance enhancements. [TR-3]

4.2.4   The Offeror should describe the capability of the system to compile and run CUDA based applications. [TR-1]

4.2.5   The system should provide MPI libraries that support MPI 4.0 or higher and make GPU-aware MPI available wherever this is supported by the GPU vendor, and this shall be capable of running a job at full-system scale. The Offeror should describe any extensions or limitations to the MPI standard in the available MPI libraries. [TR-1]

4.2.6   The system shall support the Process Management Interface (PMI). The Offeror should describe the version and supported integrations. [TR-1]

4.2.7   The Offeror should provide optimized BLAS, LAPACK, ScaLAPACK and FFT libraries for both CPUs and GPUs. Describe the offering. [TR-1]

4.2.8   The Offeror should describe any provided communication libraries (e.g., PGAS libraries and task-based programming libraries). [TR-3]

4.2.9   The Offeror should describe support for optimized scientific I/O libraries for CPUs and GPUs (e.g., HDF5, NetCDF). [TR-3]

4.2.10  The Offeror should describe any limitations of using provided libraries (used to expose high-performance CPU, GPU, I/O, or communication capabilities in 4.2.5-4.2.9) through standard Python packages (e.g., NumPy, SciPy, h5py, mpi4py) or using standard Python packaging tools. [TR-2]

4.2.11  The Offeror should describe performance optimizations for distributed or parallel execution of Python programs through Python multiprocessing, Dask, Ray, or other similar non-MPI-based runtimes. [TR-2]

4.2.12 The Offeror should describe any provided optimized libraries for execution of machine learning and AI workloads for both CPUs and GPUs, such as those required for optimal execution of deep learning frameworks like PyTorch and Tensorflow. [TR-1]

4.2.13 The Offeror should describe any support for distributed deep learning libraries (e.g., Horovod, PyTorch DDP) that enable scaling of training workloads across the full system and any provided tools that accelerate ML/AI/DL based workflows. (e.g., hyperparameter optimization, tracking experiments and integration with simulation and data pipelines). [TR-2]

4.2.14 The Offeror should describe support for direct data movement and access from the GPU to improve network bandwidth and latency (GPU-to-GPU) and I/O performance (GPU-to-Storage). [TR-3]

4.2.15 The Offeror should describe all provided profiling tools which include MPI and OpenMP profiling, and support for all user-accessible hardware (CPUs and GPUs) and any provided compilers. [TR-1]

4.2.16 The Offeror should describe support for APIs, including Linux perf, that enable profilers and other performance optimization tools to access CPU and GPU performance counters on the system. Include any restrictions on perf_event_paranoid, required kernel modules, or other security considerations. [TR-1]

4.2.17 The Offeror should describe all provided debugging tools for applications running on all user-accessible hardware such as gdb for CPUs and equivalents for GPUs. [TR-1]

4.2.18 The Offeror should provide a mechanism for users to build and run [Open Container Initiative (OCI)](#) compliant containers on the system without requiring privileged access to the system or allowing a user to escalate privilege. [TR-1]

4.2.19 The Offeror should describe how software and hardware dependencies, such as device driver libraries, MPI libraries, and libfabric, can be accessed by containers, including dynamic mechanisms to maintain accessibility of these dependencies when software updates are made. [TR-1]

4.2.20 The Offeror should describe any provided container images for users (e.g., libraries, applications), the licensing model and how they can be distributed (e.g. can we distribute a container build on top of an Offeror provided container), and the image registry where these container images may be published (ex: an internal or public registry). [TR-2]

4.2.21 The Offeror shall describe any slowdowns and scaling limitations that would be observed due to running an application in a container up to the full scale of the system. [TR-2]

## 4.3 Workflow Readiness Support

4.3.1   The Offeror should include in their proposal a separately priced plan to assist in transitioning select NERSC applications to the system (e.g., [NESAP](#) focus for simulations, data, and learning), and shall propose a vehicle (e.g. a Center of Excellence (COE)) for supporting the successful execution of this plan. Support could be provided by the Offeror and the CPU and GPU vendor. The Offeror should provide access to experts in the areas of compilers and application performance in the form of staff training and deep-dive interactions with a set of teams. Proposed plan will be used to mutually develop a **Center of Excellence for Workflow Readiness Plan** as described in Appendix B.  [TR-1]

4.3.2   The Offeror should include in their COE plan support for transitioning select workflows to the system. Support could be provided by the Offeror and/or key technology providers (e.g., the CPU and GPU vendors, storage, networking, third-party software, etc.) addressing overall workflow performance. The Offeror should include how they will collaborate with third-party developers from open source communities. [TR-2]

4.3.3   The Offeror should propose user training available during and outside of the COE for the lifetime of the system. Activities should target effective use of the user environment, performance and optimization. The description should include topics, frequency, and format (such as classroom training or online training, hackathons, etc.). Proposed training will be used to mutually develop a user **Training and Education Plan** as described in Appendix B. [TR-1]

## 4.4 Programming the Data Center

4.4.1   The system should support complex workflows described in the **NERSC-10 Workflows Whitepaper** through REST API interfaces or other mechanisms that expose functionality to users and automated services. The Offeror should describe the capabilities that their REST APIs expose, as well as how they are documented and tested.  This description may include, but is not limited to, the following capabilities: [TR-1]

- System and subsystem status and health
- Data transfer, management, and archiving
- Orchestration of workflows, persistent services, CI/CD workflows (including container deployment), and complex science workflows
- Dynamic reconfiguration of storage, compute, and networking hardware
- Any authentication and authorization requirements/expectations for their REST APIs.

4.4.2   The Offeror should describe any capabilities that enable or improve multi-tenancy support (on compute and/or WEN nodes) that goes beyond the status quo (shared node

jobs). This description may include, but is not limited to, the following capabilities:
[TR-2]

- Protecting users from one another through minimized privileges and other mitigation techniques to prevent escalations in privileges
- Virtualization and container networking (e.g., SR-IOV, VXLAN), including details of hardware offload capabilities, the number of tenants supported and guarantees of isolation between tenants.

4.4.3    The system should support integrating with external cloud capabilities provided by the Offeror and/or commercial cloud service providers to enable resilient workflows, urgent computing or specialized services (see the NERSC-10 Workflow Whitepaper). The Offeror should describe any support for hybrid-cloud capabilities, including bursting to cloud resources, portability of programming environments, data management and movement, and access to specialized hardware and services.  [TR-3]

4.4.4    The Offeror should describe capabilities to support "server-less" or Function-as-a-Service, including how it could integrate with the system and scheduler, security model, scaling and performance. Describe any capabilities to support an event-based message model that can be used to publish and subscribe to system events, job events, data-related events, and other event types that can be integrated into and used to support complex workflows. [TR-3]

4.4.5    The Offeror should describe any provided specialized hardware or integrated technologies to enhance support for composable workflows. The hardware could be present in the WEN and/or compute partitions of NERSC-10. The description should outline the programmable capabilities, the intended scope (e.g., user programmable vs a secure platform managed by trusted identities), the available memory, processing hardware, interfaces, and hardware acceleration capabilities. [TR-3]

4.4.6    The Offeror should describe data lifecycle management capabilities, such as policy-driven data movement and capability to integrate with a site-wide scheduling resource. Describe any methods for enabling persistent user-defined metadata to methods enable tracking and sharing data across the entire storage ecosystem or for automatically attaching workflow-based metadata to files.  [TR-3]

4.4.7    The Offeror should describe how their proposed system could interface with quantum computing hardware, including any software and library support. Offeror should also describe any available software they have to support quantum computing simulation at scale. [TR-3]

# 5.0 SYSTEM SOFTWARE & MANAGEMENT

5.0.1    The NERSC-10 system should include management capabilities that facilitate integration with the evolving HPC workflow-driven environment. The management system should

- employ configuration management to ensure reproducibility and automation of critical tasks (e.g., continuous deployment of operating system images, container images, microservices on compute nodes, server nodes and any support devices, and reinstallation when necessary for operational reasons).
- employ software components that should not restrict the evolution of the NERSC-10 workflow ecosystem, comply with open standards when available (e.g., Redfish) and provide documented programming interfaces. For some components, NERSC may choose to use open source or 3rd-party software over the course of the production lifecycle.

The Offeror should provide a high-level overview of their proposed system management solution and any limitations toward achieving a modular environment. [TR-1]

## 5.1 Infrastructure Services

5.1.1. The Offeror should describe remote manageability capabilities of the compute nodes, network switches, platform and QoS storage, power distribution units and servers comprising the system, including power control and console access, firmware updates, zero-touch provisioning, diagnostics, event logs, and alert capabilities. These capabilities should be accessible via documented APIs, preferably based on open standards, and a user interface. [TR-1]

5.1.2. The Offeror should describe any features provided for scalable full-platform management software that automates the management of all hardware, provides a comprehensive overview of system operations and automates whole-system maintenance actions. Relevant features include, but are not limited to, sequenced power up and power down of the system; summarization of temperature, power and other sensors; automating firmware and configuration updates; maintaining an inventory of field-replaceable units over the system lifetime; collecting alert and error information from hardware. [TR-1]

## 5.2 Operating System

5.2.1. The system should include a commercially-supported, non-proprietary Linux operating system (OS) environment on all visible service partitions (e.g., front-end nodes, service nodes, I/O nodes). The Offeror should describe the proposed Linux environment. [TR-1]

5.2.2. The system should include an optimized compute partition operating system such that all optimizations are limited to opensource linux kernel modules that can be rebuilt onsite, to provide an efficient execution environment for applications running up to full-system scale. The Offeror should describe any HPC relevant optimizations made to the compute partition operating system. [TR-1]

5.2.3. The Offer should enable all provided device drivers or kernel modules to be rebuildable and manageable by the University with the operating system proposed in 5.2.1. [TR-1]

5.2.4. The Offeror should provide access to source code, and necessary build environment, for all software except for firmware, compilers, and third-party products. The Offeror should provide updates of source code, and any necessary build environment, for all software over the life of the subcontract. [TR-1]

## 5.3 Platform Management

5.3.1. The Offeror should describe the system configuration management and diagnostic capabilities of the system that address the following details of system management: [TR-1]

- Any effect or overhead of software management tool components on the CPU or memory available on compute nodes.
- Support for multiple simultaneous or alternative system software configurations, including estimated time and effort required to install both a major and a minor system software update.
- User activity tracking, such as audit logging and process accounting.
- Unrestricted privileged access to all hardware components delivered with the system.

5.3.2. The system should have no single points of failure that would cause a **system outage** (defined in the appendix).  The system should remain in an operational or degraded state after the unexpected failure of, or planned maintenance on, any single FRU, server or switch and during any repair or other maintenance action. The Offeror should describe RAS capabilities to mitigate single points of failure (hardware or software) and the potential effect on running applications and system availability. [TR-1]

5.3.3. The Offeror should describe the resilience, reliability, and availability mechanisms and capabilities of the system to mitigate any condition or event that can potentially cause a job interrupt and how a job maintains its resource allocation and is able to relaunch an application after an interrupt.  [TR-2]

## 5.4 System Software Deployment

5.4.1 The system should include the ability to perform rolling upgrades and rollbacks on a subset of the system while at least half of the system remains in production operation. The Offeror should describe the mechanisms and limitations of the continuous deployment framework. [TR-1]

5.4.2 The Offeror should describe the process for scalable boot, reconfiguring and rebooting of compute, server and any other node types in the system. The description should include an overview of the node boot process (warmboot and coldboot), including

secure boot, stateless/stateful node provisioning, and infrastructure automation for customization and configuration of a node, the coordination, ordering and parallelism of the boot process, and techniques to provide rapid configuration and rebooting. Include how the time required to reboot scales with the number of nodes being rebooted. [TR-2]

5.4.3    The Offeror should describe any suggested system development tools to make deployments easier, for example container registry, container image management, automated testing and version control, and describe how it integrates with the system management workflow. [TR-3]

5.4.4    The Offeror should describe how NERSC can add new 3rd party hardware to the system, for example an AI accelerator partition, fabric attached memory, specialized storage, and a cluster of Linux servers. The description shall include any requirements, interfaces or standards that must be provided by this third-party hardware for addition to the system. Describe any provided specialized networks for resources requiring higher bandwidth and lower latency than the HPC network would provide.  [TR-3]

## 5.5 Data Collection and Monitoring

5.5.1    The system should include a secure mechanism whereby all monitoring data and logs captured are available to the University, and will support an open monitoring API to facilitate lossless, scalable sampling and data collection to publish and subscribe to monitoring data. Any filtering that may need to occur will be at the option of the University. The system should include a sampling and connection framework that allows the system manager to configure independent alternative parallel data streams to be directed off the system to site-configurable consumers. [TR-1]

5.5.2    The system should include mechanisms to collect, provide, store and generate alerts to monitor the status, health, and performance of the system.  These mechanisms and data should adhere to available open standards (when available), be open source or provide documented APIs and data definitions if only a proprietary solution is available. The Offeror should describe these capabilities, which should include at least the following: [TR-1]
   ● Environmental measurement capabilities for all systems and peripherals and their sub-systems and supporting infrastructure, including power, energy consumption, voltage, cooling and temperature, including sampling frequency, accuracy of the data, and timestamps of the data for individual points of measurement.
   ● Metrics related to memory, network and other error correction or faults.
   ● Metrics of both HPC protocols and TCP/IP flows.  This shall include switch and/or router data, load balancing, error counters, congestion state, throttling, throughput, and latency for select packets traversing the network(s).
   ● Resource utilization for memory, cpu, network, storage and accelerator devices.

- The system as a whole, including all levels of integrated and attached storage, and their associated hardware performance counters, degraded components and impending failure.

5.5.3  The Offeror should provide tools for the collection, analysis, integration, and visualization of metrics and logs produced by the system (e.g., peripherals, integrated and attached storage, and environmental data, including power and energy consumption). [TR-1]

# 6.0 SYSTEM NETWORKS

6.0.1  The Offeror should describe the high speed network including: [TR-1]

- High level description of how traffic would be routed and protocol translations or bridges.
- Scale of transfers and number of connections may be established per network interface and in aggregate.
- Aggregate bandwidth and transfer rates between compute nodes and the other networks.

6.0.2  The Offeror should provide a lower-level communication (LLCA) API  in the form of either UCX (https://openucx.org/) or libfabric (https://ofiwg.github.io/libfabric/).  The Offeror should describe any enhancements or limitations that can be expected in meeting full support of the standards and latest version of the LLCA. [TR-1]

6.0.3  The Offeror should describe link failure resilience throughout the network.  The number of links, network-interfaces and switch failures that can occur while maintaining connectivity and how performance degrades as links fail.  [TR-1]

6.0.4  The Offeror shall describe the out-of-band management network and mechanisms to securely extend management segments through an intermediary network, such as a data center network. [TR-1]

6.0.5  The Offeror shall provide a management platform that enables dynamic addition/removal of network segments while maintaining state tracking and recovery. [TR-3]

6.0.6  The Offeror shall describe mechanisms that provide Quality of Service for the interconnect (such as congestion control and traffic classes/virtual channels).  The Offeror shall describe how these are configured by the University and the guarantees they provide. Of particular interest are mechanisms that enhance the NERSC-10 systems performance on the complex workflows described in the **NERSC-10 Workflow Whitepaper.** [TR-2]

6.0.7  The Offeror should provide a high-bandwidth and resilient solution that allows external connectivity to and from the system to the NERSC data center Ethernet

network with support for thousands of simultaneous transfers and is capable of at least 26 Terabits per second aggregate throughput. NERSC will provide file systems configured for different purposes to users including, but not limited to, a HOME file system. The Offer should work with the University to ensure the NERSC-10 system mounts these file systems to achieve a high level of user satisfaction.  [TR-1]

6.0.8   The Offeror should describe proposed support for: [TR-1]

- Jumbo frames, IPv6, IPv4, TCP/IP, UDP and virtual networks support;
- the ability to control IP traffic using Access Control Lists;
- the ability to load balance and route traffic across multiple paths, and
- the ability to dynamically configure routing and exchange routing information between the data center, storage and other external networks.

6.0.9   The Offeror should describe how the external connectivity solution will utilize IETF standards-compliant technology for functionality that includes (but is not limited to): [TR-2]
- Congestion control mechanisms across network boundaries
- QoS required for reliable connections, with configurable buffers
- Traffic draining capability at a link or adjacency level, with flow tracking ability using sampled flow or similar
- Encryption and authentication
- Integrated SDN with job management and scheduling systems.

6.0.10  The Offeror should describe support for CXL 2.0 (or greater) standard in its CPU-, GPU-, NIC & storage offering. Include the level of support and the intended usage models of CXL.io, CXL.mem and CXL.cache.  [TR-3]

# 7.0 STORAGE SYSTEMS

The NERSC-10 system requires two distinct storage system offerings to accommodate diversity of I/O needs.

The platform storage system (PSS) should accommodate the I/O and storage needs of multiple workload types. The PSS should be designed for scale, capable of serving I/O to whole-system jobs (as well as smaller jobs), maximize bandwidth and have a data protection scheme in place (e.g., a form of RAID or erasure coding).

The QoS storage system (QSS) should be designed to provide deterministic performance through the use of a configurable Quality of Service (QoS) mechanism. The QSS will be designed for scale, capable of serving I/O to whole-system jobs as well as smaller jobs. Of particular interest are mechanisms that seamlessly enable the time-sensitive workflows described in the **NERSC-10 Workflow Whitepaper**.

Both PSS and QSS will also be accessed via other systems, e.g. external data transfer nodes, and as such, should not require the compute system for them to be accessible, and vice versa.

## 7.1 Platform Storage System (PSS)

7.1.1     The Offeror should propose a separately priced Platform Storage System solution at least 120 PB in size with an aggregate bandwidth of at least 20 TeraBytes/s. [TR-1]

7.1.2     The Offeror should describe the scalable unit for capacity, bandwidth, IOPS, and provide details for increasing each, and any associated limitations, including number of concurrent clients. Describe projected characteristics of primary storage devices such as media type, usable capacity, storage interfaces (e.g., NVMe, PCIe), and media durability. [TR-1]

7.1.3     The Offeror should describe all available interfaces to platform storage for the system, including but not limited to POSIX, Kubernetes CSI, and other APIs. Describe any exceptions to POSIX compliance, time to consistency, any potential delays for reliable data consumption. [TR-1]

7.1.4     The Offeror should propose a method(s) for PSS to be accessible even if the compute system is unavailable. Describe any scenarios where a rebalance of resources is required after a reconfiguration and any scenarios where downtime for the entire PSS is required. [TR-1]

7.1.5     The Offeror should describe data protection schemes, mechanisms for recovery and related performance impact. Describe any anticipated loss of performance over time as the file system ages or reaches capacity. [TR-1]

7.1.6     The Offeror should provide features to enforce and report upon soft (accounting) and hard (enforcement) quotas based on uid, gid or other constructs. [TR-1]

7.1.7     The Offeror should provide system features for metadata scanning, metadata processing and file/object purging across the entire PSS to allow for purging of older data and provide data to feed user data management tools. Describe the expected rate and elapsed time of a full-system scan and any performance impact while metadata scan or purge is happening. [TR-1]

7.1.8     The Offeror should describe the ability of the PSS to purge based upon individual files' last access time, modification time, owner, group, location and size including capability to explicitly include or exclude files and directories, and a dry-run reporting mode. [TR-2]

## 7.2 QoS Storage System (QSS)

7.2.1     The Offeror should propose a separately priced QoS Storage Solution, at least 80 PB in size, which shall provide storage and deterministic I/O to the system, to provide consistent performance and ensure timely results.

Two modes of operation are anticipated:

- Node-local disk emulation by assigning sufficient capacity and performance to a mounted directory (container, virtual disk).
- Parallel I/O and/or shared work directories by assigning capacity and performance to an allocation shared across nodes.

Successful implementation may require coordinating the QSS with network QoS to guarantee end-to-end performance and minimize the effects of contention from other jobs. Describe how the design goals will be achieved. [TR-1]

7.2.2   The Offeror should describe the scalable unit for capacity, bandwidth, IOPS, and provide details for increasing each, and any associated limitations, including number of concurrent clients. Describe projected characteristics of primary storage devices such as media type, usable capacity, storage interfaces (e.g., NVMe, PCIe), and media durability. [TR-1]

7.2.3   The QSS shall provide user and administrator control mechanisms for creating, serving, and querying the status of storage allocations, where an allocation is a directory (container, virtual disk, …) that has a capacity, capabilities (e.g., minimums and maximums for bandwidth and IOPS), and policies (e.g., persistence or duration, user quota, namespace mapping).  Additionally, the QSS shall provide a mechanism to display current unallocated capacity and capability.  [TR-1]

7.2.4   The Offeror should describe how allocation mechanisms could be automated, or otherwise improved in the future. For example, how the QSS could receive and serve requests for storage allocations via Slurm. [TR-3]

7.2.5   The Offeror should describe all available interfaces (and any associated limitations) to QSS, including but not limited to POSIX, S3, Kubernetes CSI, NVMEoF, and other APIs. Describe any exceptions to POSIX compliance, time to consistency, and any potential delays for reliable data consumption. [TR-1]

7.2.6   The Offeror should describe any data protection schemes, mechanisms for recovery, and related performance impact. Describe any anticipated loss of performance over time as the file system ages or reaches capacity. [TR-1]

7.2.7   The Offeror should propose a method(s) for QSS to be accessible even if the compute system is unavailable. Describe any scenarios where a rebalance of resources is required after a reconfiguration and any scenarios where downtime for the entire QSS is required. [TR-1]

7.2.8   The Offeror should provide system features for metadata scanning and processing across the entire QSS to produce data to feed tools for user data analysis and management. Describe the expected rate and elapsed time of a full-system scan and any performance impact while metadata scan is happening. [TR-2]

# 8.0 SYSTEM OPERATION

## 8.1 Resilience, Reliability and Availability Metrics

The ability to achieve the NERSC-10 mission goals hinges on the productivity of system users. System availability is therefore essential and requires system-wide focus to achieve a resilient, reliable, and available system. **For each metric specified below, the Offeror should describe how they arrived at their estimate(s).**

8.1.1   The system is available if it meets the system outage criteria in the glossary. The Offeror should propose a system availability. [TR-1]

8.1.2   The Offeror should propose a System Mean Time Between Interrupt (SMTBI). [TR-1]

8.1.3   The Offeror should propose a Job Mean Time To Interrupt (JMTTI) for a single job running on the entire system. [TR-1]

8.1.4   The system should complete power on and power off in a timely manner. The Offeror should describe the sequence of steps and timings for full system initialization and full system shutdown. Include any dependencies and how timings may scale with the size of the system. [TR-1]

## 8.2 System Security

8.2.1   The Offeror should describe the security capabilities of the proposed compute node and service partition operating systems.  [TR-1]

8.2.2   The Offeror should describe how the system may be configured to support Zero-Trust requirements as described in the CISA Zero Trust Maturity Model (https://www.cisa.gov/zero-trust-maturity-model) [TR-1]

8.2.3   The Offeror should describe how their implementation of Identity Management (including Federated Identity) for the system works, what protocols and standards are being utilized (e.g., SpiFFE, SPIRE) and what system services and/or service accounts will use that framework. [TR-1]

8.2.4   The Offeror should describe how they will implement a Root of Trust to assure a boot environment for the system that is secure. [TR-1]

8.2.5   Security vulnerabilities in the software supplied by the vendor should be addressed with patches and/or updates in a timely manner depending on the classification and severity of the vulnerability. The Offeror should describe the process for handling security vulnerabilities and the time to provide patches from confirmed availability for: [TR-1]

- non-critical vulnerabilities,
- vulnerabilities as defined by CISA, and

- Common Vulnerabilities and Exposures (CVEs) in the National Vulnerability Database (NVD) as defined by NIST with a score of Critical or High as defined by the latest version of CVSS.

8.2.6 The Offeror should describe how the system validates its software and data components in compliance with the Presidential Executive Order on Software Bill of Materials.  (Cybersecurity Executive Order 14028.) [TR-1]

8.2.7 The Offeror should describe how the baseline configuration is tested, validated and documented. Describe the process of re-validation, at any time, of a running system against the baseline, including how the system can be audited so its current state can be captured and documented. [TR-2]

8.2.8 The Offeror should describe the security model, tools and functionality for observability for any system applications and services that are containerized. [TR-1]

8.2.9 The Offeror should describe the available security controls and how they would be deployed and used for access methods to the system and/or services (for example, interactive access using ssh, non-interactive methods using APIs or Web interfaces) [TR-1]

## 8.3 Power and Energy

8.3.1 The maximum power consumed by the system and its peripheral systems, including the proposed storage systems, will not exceed 20 MW. The maximum power consumption includes all equipment provided by the proposal. The Offeror should describe how its proposal fits within the power budget. [TR-1]

8.3.2 The Offer should describe the power management capabilities of the system available to users or administrators. The description should include a description of all system control capabilities to affect power or energy use (system, rack/cabinet, board, node, component, and sub-component level) and the reaction time of this mechanism. [TR-1]

8.3.3 The Offeror shall describe the overall system measurement capabilities that enable meeting Green 500 requirements at Level 2 or higher for power, current, and voltage. [TR-3]

## 8.4 Maintenance and Support (Hardware/Software)

8.4.1 The Offeror should propose and separately price maintenance and support for all systems for a period of four (4) years from the date of acceptance of the system  in coordination with the University and subject to review. The maintenance and support will include all features outlined in the **Key Elements of the Maintenance and Support Plan** described in Appendix B or describe how the proposed plan differs. [TR-1]

8.4.2 The system should include a means for tracking and analyzing all software updates, software and hardware failures, and hardware replacements over the lifetime of the system. [TR-2]

8.4.3 Offeror should also propose additional maintenance and support extension for years 5-7. [TO-1]

8.4.4 The Offeror should provide one (1) Systems Operations and Advanced Administration training for each system delivered at facilities specified by the University. The Offeror should describe additional training available for systems operations and advanced administration available for the lifetime of the system, including topics, duration, and proposed timing. [TR-1]

## 8.5 Documentation

8.5.1 The Offeror should describe documentation provided to effectively administer and use the system, including types of documentation, format (such as user manuals, man pages, release notes, stable urls, plain text vs pdf, vendor websites, any interactive elements, etc.), initial delivery, and frequency of updates, for the following: [TR-1]

- documentation for each delivered system describing the configuration, interconnect topology, labeling schema, hardware layout, etc. of the system as deployed before the commencement of system acceptance testing

- documentation of the proposed solution to the operators and system administrators to effectively operate and configure the platform

- Documentation for users describing the programming environment and software tools, including compatibility across system software updates.

8.5.2 The Offeror should grant the University use and distribution rights for provided documentation, training session materials and recorded media to be shared with DOE Lab staff and all authorized users and NERSC support staff. The University may, at their option, make audio and video recordings of presentations from Offeror at public events targeted at the NERSC user communities (e.g., user training events, collaborative application events, Hackathons, Best Practices discussions) and make them available to all NERSC users. [TR-2]

8.5.3 All documentation shall be distributed and updated electronically and in a timely manner that maintains the productivity and performance of the system. For example, changes to the system should be accompanied by relevant documentation, such as binary compatibility after major OS upgrades. [TR-1]

8.5.4 Documentation of changes and fixes may be distributed electronically in the form of release notes. Reference manuals may be updated later, but effort should be made to keep all documentation current. [TR-1]

# 9.0 FACILITIES AND SITE INTEGRATION

The following section addresses the facility-based requirements for the proposed system. It includes pertinent information and vendor requirements for the physical, electrical, cooling, seismic, safety, and transportation aspects of designing, delivering, installing, and integrating the system at the facility.

9.0.1   The Offeror's proposed system and integration plan will include, but is not limited to, all features outlined in the **Site Preparation Facilities Plan** in Appendix B and should conform to the **Facility and Site Integration Specifications** provided in Appendix C. The Offeror should describe any limitations to meeting the specifications. The successful offeror should work with the University to provide site integration information in a timely manner.  [TR-1]

9.0.2   The Offeror should describe capabilities to improve the environmental sustainability and system efficiency, beyond energy efficiency, through the entire lifecycle of the system including design, manufacturing, deployment, and operation, and the ability to reuse and recycle components into future systems, and in final disposal. [TR-2]

9.0.3   The Offeror should propose to deinstall, remove and/or recycle the system and supporting infrastructure at end of life.  Storage media shall be wiped or destroyed to the satisfaction of the University, and/or returned to the University upon request. [TO-1]

# 10.0 DELIVERY AND ACCEPTANCE

The Delivery and Acceptance Test Plan will describe the steps needed to validate the system, including pre-delivery, post-delivery, and acceptance test plan for the NERSC-10 Production system and supporting systems, such as Early Access System (EAS) and Test and Development System (TDS).

Acceptance testing may comprise multiple components for which the overall goal is to ensure that the system as a whole is high performance, scalable, resilient, and reliable. Acceptance testing may exercise the system infrastructure with a combination of functionality tests, performance tests, forced failures, and stability tests. Any requirement described in the Technical Specification may generate a corresponding acceptance test element. The specifics of an acceptance plan will be determined before system delivery as part of the NERSC-10 Acceptance Test Plan Milestone.

10.0.1  The Offeror should work with the University to define the pre-delivery, post-delivery and acceptance test plan as part of the NERSC-10 System Acceptance Test Plan Milestone. A Sample Acceptance Test is provided in Appendix A. [TR-1]

10.0.2  The NERSC team and the Successful Offeror should perform pre-delivery testing at the factory on the hardware to be delivered. Any limitations for performing the pre-delivery testing should be identified in the Offeror's proposal, including scale and

licensing limitations (if any). During pre-delivery testing, the Successful Offeror should: [TR-1]

- Demonstrate RAS capabilities and robustness using simple fault injection techniques, such as disconnecting cables, powering down subsystems, or installing known bad parts.
- Demonstrate functional capabilities on each segment of the system built, including the capacity to build applications, schedule jobs, and run them using a customer-provided testing framework. The root cause of application failure must be identified prior to system shipping.
- Provide a file system sufficiently provisioned to support the suite of tests.
- Provide onsite and remote access to the NERSC team to monitor testing and analyze results.
- Instill confidence in the ability to conform to the statement of work

10.0.3 The NERSC team and the Successful Offeror staff shall perform site integration and post-delivery testing on the fully delivered system. [TR-1]

- During post-delivery testing, the pre-delivery tests shall be run on the full system installation.
- Where applicable, tests shall be run at full scale.

10.0.4 The NERSC team and the Successful Offeror staff shall perform onsite acceptance testing on the fully installed system. [TR-1]

# 11.0 PROJECT AND RISK MANAGEMENT

The development, pre-shipment testing, installation and acceptance testing of the NERSC-10 System and the management of the Non-Recurring Engineering (NRE) subcontract are complex endeavors and will require close cooperation between the successful Offeror and the Laboratory. The documents described in this section are not required in the RFP response, however a commitment from the Offeror to deliver the documents within the timeframe described required.

11.0.1 The Offeror should commit to develop, maintain and submit to the University the Planning Deliverables described in the **Key Project Planning Deliverables** section and form the relevant **Project Working Groups**. [TR-1]

11.0.2 The Offeror should provide in its RFP response a set of milestones in this section as described in the **Key Milestone Dates.** [TR-1]

11.0.3 LBNL and Offeror should schedule and complete a **Project Planning Kickoff Meeting** to mutually understand and agree upon project management goals, techniques, and processes for the NERSC-10 system and the NRE subcontract. The kickoff meeting shall take place no later than 45 days after contract award. [TR-1]

**DESCRIPTIONS:**

**Key Project Planning Deliverables**

The Offeror will develop, deliver, submit for approval and maintain the following Planning Deliverables. These plans are described in more detail in Appendix B: Project and Risk Management - Key Planning Elements Descriptions. Initial versions and updates of these plans shall be provided in specific agreed upon time frames. Each of the plans and any revisions will be submitted for comment and approval to the University's project leadership.

- Project Plan
- Communication Plan
- Risk Management Plan
- Site Preparation Facilities Plan
- Chemical Management Plan
- Network Plan
- Acceptance Test Plan
- System Delivery and Installation Plan
- Training and Education Plan
- Center of Excellence for Workflow Readiness Plan
- Maintenance and Support Plan

**Project Working Groups**

As described above, the NERSC-10 contracts must represent a partnership that is committed to delivering the most useful system possible. Upon subcontract award, the selected Offeror and the University will assess the project for areas in which deep collaboration is necessary to ensure meeting that goal. The partnership will form working groups (WGs) for these topics. Each WG will interact in regard to all details of the technical topic; the selected Offeror will not attempt to limit the scope of these interactions. Specific details include NRE deliverables related to the technical topic and deployed software and hardware in the systems being built. The WGs will serve as a key conduit to identify, to refine and to understand NERSC requirements in detail and to ensure that the delivered system meets those requirements to the greatest extent possible. WGs will establish a regular schedule for electronic meetings (e.g., telecons). Each Quarterly face-to-face meeting may include WG breakout sessions. WG breakout sessions will be determined by project management, including the University and selected Offeror representatives. Project management will regularly assess WG progress and identify topics for which WGs are no longer required or additional topics for which new WGs are needed.

**Key Milestone Dates**

Offeror will provide the University, in its proposed response, a proposed set of milestones for this section and, for each milestone, a proposed associated payment that is applicable to Offeror's proposed development and deployment timeline and methodology. Offeror is encouraged to identify milestones for each year of the project that merit revenue that the Offeror can legally recognize in that year.

The successful Offeror and University will hold a Technical Decision Point evaluation and joint planning meeting 9-12 months before System delivery. At the Technical Decision Point meeting the final configuration of the System will be determined based on technology status and evaluations and component pricing. Performance targets will be re-evaluated and converted into requirements.

Prior to award, the Offeror and University will finalize the list of Key Milestone Dates, including dates for necessary Technical Decision Point evaluations. Following is a list of the kinds of key dates of importance to the University. Other key dates may be needed for phased installations or deployments featuring major upgrades during the subcontract. <u>Early completion is highly desired</u>.

- Project Liaisons assigned to include the Offeror Project Manager, Executive Point of Contact, Service Manager, Contract Manager and Account Manager
- Project Plan complete;
- Pilot system delivery;
- Technical Decision Point to exercise any proposed system options;
- IO Subsystem late binding decision(s) to exercise any proposed tier options;
- System Delivery and Installation Plan to include date of on-site support personnel to arrive on site, e.g., hardware, storage and software specialists;
- Begin delivery and installation of system and exercised storage and network options;
- System Installation and Integration complete including IO Subsystem;
- System Accepted including IO Subsystem.

# APPENDIX A: Sample Acceptance Test Plan

**A.1 Staffing and Safety**

The Subcontractor shall provide sufficient staffing to perform the installation, initial testing and acceptance of the System.

The Subcontractor shall conform to the University's safety protocols and policies and complete all necessary documents (e.g,. approved safety plan) and required training. The system should be installed in accordance with the University's safety policies and the approved joint safety plan. The Subcontractor shall designate a safety supervisor who will monitor Subcontractor staff for adherence to the University's safety policies and the approved safety plan.

**A.2 Pre-Delivery Assembly, Quality Assurance and Factory Test Plan**

The Subcontractor shall perform the pre-delivery assembly and quality assurance tests of the System and agreed-upon sub-configurations at the Subcontractor's location demonstrating all hardware is fully functional prior to shipment. A factory test plan will be agreed on by the University and Subcontractor 30 days prior to the factory test.

At its option, the University may send a representative(s) to observe testing at the Subcontractor's facility. Work to be performed by the Subcontractor includes:

- All hardware installation and assembly
- Burn in of all components
- Installation of software
- Implementation of the University-specific production system-configuration and programming environment necessary to complete required testing
- Perform tests and benchmarks to validate functionality, performance, reliability, and quality
- Run benchmarks and demonstrate that benchmarks meet performance commitments

The pre-delivery test shall consist of (but is not limited to) the following tests:

| Name of Test | Pass Criteria |
|---|---|
| System power up | All nodes boot successfully |
| System power down | All nodes shut down |
| Monitoring | Monitoring software shows status for all nodes |
| Reset | "Reset" functions on all nodes |
| Power On/Off | Power cycle all components of the entire system from the console |
| Fail Over/Resilience | Demonstrate operation of all fail-over or resilience mechanisms |

| Benchmarks | The system should demonstrate the ability to achieve the required performance level on all benchmark requirements. |
|---|---|
| 72 Hour test | High availability of the production system for a 72-hour test period under constant throughput load |

### A.3 Site Integration and Post-Delivery Testing

The system should be delivered, installed, fully integrated, and shall undergo Subcontractor stabilization processes. Post-delivery testing shall include replication of all of the pre-delivery testing steps, along with appropriate tests at scale, on the fully integrated platform.

When the Subcontractor has declared the System to be stable, the Subcontractor shall make the System available to University personnel for site-specific integration and customization. After the Subcontractor's System has completed site-specific integration, security screening and customization, the acceptance test shall commence.

### A.4 Acceptance Testing

The Acceptance Test Period shall commence when the System has been delivered, physically installed, undergone stabilization and site-specific integration and customization, and conforms to all requirements in the statement of work designated "critical" priority. The Acceptance Test Period should target 60 days or until it has met agreed upon exit criteria.

The Subcontractor shall not be responsible for failures to meet the performance metrics or the availability metrics set forth in this Section, if such failure is the direct result of modifications made by the University to Subcontractor source code. Such suspension will be only for those requirements that fail due to the modification(s) and only for the length of time the modification(s) result(s) in the failure.

The Subcontractor shall supply source code used, compile scripts, output, and verification files for all tests. All such provided materials become the property of The University.

All tests shall be performed on the production configuration of the System, as it will be deployed to the University user community. The University may run all or any portion of these tests at any time on the System to ensure the Subcontractor's compliance with the requirements set forth in this document.

The acceptance test shall consist of Functionality Demonstration, System Test, System Resilience Test, Performance Test, and an Availability Test, performed in that order.

**Functionality Demonstration.** The Subcontractor and the University will perform the Functionality Demonstration on a dedicated system. The Functionality Demonstration shall show that the System is configured and functions in accordance with the statement of work. Demonstrations shall include, but are not limited to, the following:

- Remote monitoring, power control and boot capability

- Network connectivity
- File system functionality
- Batch system
- System management software
- Program building and debugging (e.g., compilers, linkers, libraries, etc.)

**System Test.** The Subcontractor and the University will perform the System Test on a dedicated system. The System Test shall demonstrate that the System is configured and functions in accordance with the statement of work. Demonstrations shall include, but are not limited to, the following:

- Two successful System cold boots to production state in accordance with required timings, with no intervention to bring the System up. Production state is defined as running all System services required for production use and being able to compile and run parallel jobs on the full System. In a cold boot, all elements of the System (compute, login, I/O, network) are completely powered off before the boot sequence is initiated. All components are then powered on.
- Single node power-fail/reset test: Failure or reset of a single compute node shall not cause a system-wide failure. A node shall reboot to production state after reset in accordance with required timings.

**System Resilience Test.** The Subcontractor and the University will perform the System Resilience Test on a dedicated System. The System Resilience Test shall show that the System is configured and functions in accordance with the statement of work.

All demonstrable System resilience features of the system should be demonstrated via fault-injection tests when running test applications at scale. Fault injection operations should include both graceful and hard shutdowns of components. The metrics for resilience operations include correct operation, any loss of access or data, and time to complete the initial recovery plus any time required to restore (fail-back) a normal operating mode for the failed components.

**Performance Test.** The Subcontractor and the University shall run the NERSC-10 tests and workflow component benchmarks, file system tests, a minimum of five times each, and as described in the Benchmark Run Rules. Benchmark answers must be correct, and each benchmark result must meet or exceed performance commitments.

Benchmarks must be run using the supplied resource management and scheduling software. Except as required by the run rules, benchmarks need not be run concurrently. If requested by the University, Subcontractor shall reconfigure the resource management software to utilize only a subset of compute nodes, specified by the University. Performance must be consistent from run to run.

**Availability Test.** The Availability Test will commence after successful completion of the Functionality Demonstration, System Test and Performance Test. The Subcontractor shall

perform the Availability Test; at this time or before, the University will add user accounts to the System. The Availability Test shall be 30 contiguous days in a sliding window within the Acceptance Test Period. The NERSC-10 System must demonstrate the required availability of the System.

During the Availability Test, the University shall have full access to the System and shall monitor the System. The University and users designated by the University shall submit jobs through the NERSC-10 resource management system. These jobs shall be a mixture of benchmarks from the Performance Test and other applications.

The Subcontractor shall adhere to the System Availability and Reliability requirements as defined below:

- All hardware and software shall be made fully functional before the availability test can be declared complete. Any down time required to repair failed hardware or software shall be considered an outage unless it can be repaired without impacting system availability.
- Hardware and software upgrades shall not be permitted during the last 7 days of the Availability Test. The system should be considered down for the time required to perform any upgrades, including rolling patch upgrades.
- No critical bugs shall be open during the last 7 days prior to the Availability Test.
- During the Availability Testing period, if any System software upgrade or significant hardware repairs are applied, the Subcontractor shall be required to run the Benchmark Tests and demonstrate that the changes incur no loss of performance. At its option, the University may also run any test or benchmark deemed necessary. Time taken to run the Benchmark and other tests shall not count as downtime, provided that all tests perform to specifications.
- Every test in the Functionality Test, Performance Test and NERSC-defined workload shall obtain a correct result in both dedicated and non-dedicated modes.
- Each benchmark in the Performance Test shall meet or exceed the performance commitment and variation requirement.

# APPENDIX B: Project and Risk Management - Key Planning Elements Descriptions

Each key planning element is described in detail below. The specific details are designed to help the successful Offeror successfully meet its commitment, to help the University track the NERSC-10 project, and to help the University and the selected Subcontractor to understand and to mitigate risks successfully.

**<u>Key Elements of the Project Plan</u>**

The Offeror and University shall have a joint project planning call no more than 30 days after subcontract award. The Offeror should provide the University with a detailed Project Plan no later than 90 days after subcontract award that addresses, at a minimum, the following:

- The Offeror should appoint a Program Manager (PM) for the purposes of executing the Project Plan for the system and NRE contract on behalf of the Offeror. The Program Manager shall serve as the primary interface for the University, managing all aspects of the Subcontract in response to the program requirements.

- Project Management Organization Breakdown Structure (OBS) with management team's roles and responsibilities clearly defined

- Points of Contacts to include the Offeror's Technical Contact(s), Service Manager, Contract Manager and Account Manager

- Work Breakdown Structure (WBS) to include all major subsystems, each software product and each major equipment deliverable to the University

- Full Project Schedule Gantt chart for the duration of the contract

The Project Plan should be updated as timelines for delivery and installation become firm.

**<u>Key Elements of the Communication Plan</u>**

The project planning kick-off meeting shall take place no later than 45 days after contract award. A Communication Plan shall be developed and shall describe the types of communications, meetings, and progress reviews as described below:

- Daily Communication (System Contract)

  The Offeror's PM (or designate) is the owner of this meeting with a target duration one-half hour. Both Offeror and the University may submit agenda items for this meeting. These daily communications shall commence 30 days before expected system delivery and continue until both parties agree they are no longer needed

- Weekly Status Meeting (System and NRE contract)

  The Offeror's PM shall schedule this meeting with a target duration of one hour. Attendees normally include the Offeror's PM, Service Manager, University's Procurement Representative, Technical Representative and System Administrator(s) as well as other invitees.

- Quarterly Business Reviews (System and NRE contract)

    The Offeror's PM shall schedule this meeting with a target duration of no less than six hours. Attendees normally include: Offeror's PM, Offeror's Senior Management, University's Procurement Representative, Technical Representative, selected Management, selected Technical Staff and other invitees. Topics covered will cover both NRE and System contract issues that will include:

    - Program status (Offeror to present)
    - University satisfaction (University to present)
    - Partnership issues and opportunities (joint discussion)
    - Future hardware and software product plans and potential impacts for the University
    - Participation by Offeror's suppliers as appropriate
    - Other topics as appropriate
    - Both Offeror and the University may submit agenda items for this meeting.

**Key Elements of Risk Management (System and NRE contract)**

The Offeror should provide the University with a Risk Management Plan (RMP) for the technology, schedule and business risks of the NERSC-10 project 30 days after award of the Subcontract. The RMP describes the Subcontractor's approach to managing NERSC-10 project risks by identifying, analyzing, mitigating, contingency planning, tracking, and ultimately retiring project risks. The Plan shall address both the System and the NRE portions of the project. The purpose of this RMP, as detailed below, is to document, assess and manage Subcontract's risks affecting the NERSC-10 project:

- Document procedures and methodology for identifying and analyzing known risks to the NERSC-10 project along with tactics and strategies to mitigate those risks.
- Serve as a basis for identifying alternatives to achieving cost, schedule, and performance goals.
- Assist in making informed decisions by providing risk-related information.

The RMP shall include, but is not limited to, the following components: management, hardware, software; risk assessment, mitigation and contingency plan(s) (fallback strategies). Once the plan is approved by the University, the University shall review the Offeror's RMP annually.

The Offeror should also maintain a formal Risk Register (RR) documenting all individual risk elements that may affect the successful completion of the NERSC-10 project (both System and NRE contract). The RR is a database managed using an application and format approved by the University. The initial RR is due 30 days after award of the Subcontract. The RR shall be updated at least monthly, and before any Critical Decision (CD) reviews with DOE. After acceptance, the RR shall be updated quarterly. Items in the RR include, at a minimum, mitigation strategies, impact to the NERSC-10 project, severity rating, and probability of the risk occurring.

**Key Elements of Site Preparation Facilities Plan**
Site planning shall be conducted by the Offeror's Site Engineering department. Site planning consists of the exchange of system specification documents, site floor plans, and is followed by a physical inspection of facilities. The plan shall include the floor plan and the Machine Unit Specification (MUS). The Offeror should provide Preliminary Site Preparation Facilities Plan at least one year prior to the delivery of the first equipment. At least 9 months prior to the first equipment delivery, the Final Site Preparation Plan will be delivered to the Laboratory for approval. Details of the Facility requirements are in Appendix C. Items in the plan, at a minimum, include:

- Cabinet dimensions, diagrams and cabinet weights
- Electrical requirements required by the Offeror
- System layout and cabling requirements
- Raised floor requirements and cutouts
- Cable tray requirements
- Environmental requirements
- Cooling water requirements
- Safety requirements

**Key Elements of the Chemical Management Plan**
The Offeror should develop and document a chemical management plan covering all aspects of transport, storage, filling, draining, and disposal of the chemicals, MSDS safety compliance, and required training and personal protective equipment needed for compliance with safety regulations, as required. The plan shall be documented and submitted for University review and concurrence prior to the delivery and deployment of the system.

**Key Elements of the Network Plan**
No less than 6 months prior to installation of the System, the Offeror should provide an initial draft of the system network configurations for the University to review, including all network types provided, and showing compute-to-compute, compute-to-storage, and system-to-external components connectivity.

No less than 4 months prior to installation, the University and the Offeror should finalize the network design and the Offeror should provide an up-to-date copy of the network configurations reflecting the expected-at-installation design.

**Key Elements of the Acceptance Test Plan**
The University and the Offeror will create a detailed Acceptance Test Plan one year prior to the first equipment delivery and will be updated as necessary.

Items in the plan should, at a minimum, include:

- Pre-Delivery Assembly, Quality Assurance, and Preliminary Factory Testing. The plans will include how the Offeror qualifies their vendors, factory burn in and validation test plans and a pre-ship test plan for the NERSC-10 system.
- Acceptance Testing. The plans shall consist of Functionality Demonstrations, System Tests, System Resilience Tests, Performance Tests, and an Availability Test, performed in that order.
- EAS and TDS Testing. Details for testing systems in support of the NERSC-10 production system.

An example Acceptance Test Plan is included in Appendix A.

## Key Elements of System Delivery and Installation Plan

The Offeror should provide a Preliminary Delivery and Installation Plan to the University one year prior to the first equipment delivery and will be updated as necessary. The Offeror will provide a Final Installation Plan no less than 90 days before delivery. The plan shall include:

- Core installation team and staffing plan
- Detailed delivery and installation schedule
- Equipment layout and installation sequence (multi-stage deliveries)
- Detailed integration and test plan addressing all equipment and software included in the delivery
- Safety documents required by LBNL safety processes
- Diagrams showing the internal layout of cabinets

The University shall review the plan and work with the Subcontractor to promptly resolve any issues or clarifications.

## Key Elements of Maintenance and Support Plan

The Offeror should provide a maintenance, on-site support and services plan for the life of the subcontract and shall include the following features:

- **Maintenance and Support Period**
  The Offeror should propose all maintenance and support for a period of four (4) years from the date of acceptance of the system. Warranty shall be included in the 4 years. For example, if the system is accepted on April 1, 2026 and the Warranty is for one year, then the Warranty ends on March 30, 2027, and the maintenance period begins April 1, 2027 and ends on March 30, 2030. Offeror should also propose additional maintenance and support extension for years 5-7.

- **Maintenance and Support Solution**
  The Offeror should propose a maintenance and support solution with full hardware and software support for all Offeror provided hardware components and software. The principal period of maintenance (PPM) shall be for 24 hours by 7 days a week with a four hour response to any request for service. The Offeror should provide/enable access to direct communication between NERSC staff and GPU vendor technical staff.

- **Concurrent Maintenance Techniques**
  The Offeror should use continuous operations maintenance techniques (e.g. warm swap) that avoid service disruptions. Continuous operations comprise both hardware (including servicing node hardware, cabinet hardware), and software upgrades to systems management nodes, workflow nodes, storage, and compute nodes.  These actions shall not be deemed to cause a system outage if performed with the concurrence of the University and completed in a timely manner. Six hours are permitted for cabinet-level repairs and two hours for all other repairs performed concurrently, node downtime due to concurrent maintenance is counted in calculating System availability.

- **General Service Provisions**
  The Offeror should be responsible for repair or replacement of any failing hardware component that it supplies and correction of defects in software that it provides as part of the system. At its sole discretion, NERSC may request advance replacement of components which show a pattern of failures which reasonably indicates that future failures may occur in excess of reliability targets, or for which there is a systemic problem that prevents effective use of the system. Hardware failures due to environmental changes in facility power and cooling systems which can be reasonably anticipated (such as brown-outs, voltage-spikes or cooling system failures) are the responsibility of the Offeror.

  When a component has failed in service, the Offeror should replace the component with a newly manufactured or remanufactured/fully-tested component. The Offeror should not place a component back into the main system in order to determine if a failure is transient. With University's concurrence, the Offeror may use the test system to test components.

- **Software and Firmware Update Service**
  The Offeror should provide an update service for all software and firmware provided for the duration of the Warranty plus Maintenance period. This shall include new releases of software/firmware and software/firmware patches as required for normal use. The Successful Offeror should integrate software fixes, revisions or upgraded versions in supplied software, including community software (e.g., Linux or Lustre), and make them available to NERSC within twelve (12) months of their general availability.  The Offeror should provide prompt availability of patches for cybersecurity defects.

- **Call Service**
  The Offeror should provide contact information for technical personnel with knowledge of the proposed equipment and software.  These personnel shall be available for consultation by telephone and electronic mail with NERSC personnel. In the case of degraded performance, the Offeror's services shall be made readily available to develop strategies for improving performance, i.e., patches, workarounds.

- **Problem Escalation**
  The Offeror should document severity classifications and response for hardware and software problems. The description should include the technical problem escalation mechanism based either on time or the need for more technical support in the event issues are not being addressed to the University's satisfaction. Problem escalation procedures are the same for hardware and software problems. Problems should be searchable in a database and made accessible via a web interface or for download in a standard format (e.g., csv). This capability shall be made available to all individual University staff members designated by the University.

- **On-site Parts Cache**
  The Offeror should maintain a parts cache on-site at NERSC. The parts cache shall be sized and provisioned sufficiently to support all normal repair actions for two weeks without the need for parts refresh. The initial sizing and provisioning of the cache shall be based on Offeror's Mean Time Between Failure (MTBF) estimates for each FRU and each rack, and scaled based on the number of FRU's and racks delivered. The parts cache configuration will be periodically reviewed for quantities needed to satisfy this requirement, and adjusted, if necessary, based on observed FRU or node failure rates. The parts cache will be resized, at the Offeror's expense, should the on-site parts cache prove to be insufficient to sustain the actually observed FRU or node failure rates.

- **On-Site Node Cache**
  The Offeror should also maintain an on-site spare node inventory of at least 1% of the total nodes in all of the system.   These nodes shall be maintained and tested for hardware integrity and functionality utilizing the Hardware Support Cluster defined below if provided.

- **Hardware Support Cluster**
  A test and development system (TDS) described in Section 2.6 will be used as a hardware support cluster (HSC). The HSC shall support the hot spare nodes and provide functions such as hardware burn-in, problem diagnosis, etc. The Offeror should supply sufficient racks, interconnect, networking, storage equipment and any associated hardware/software necessary to make the HSC a stand-alone system capable of running diagnostics on individual or clusters of HSC nodes.

**Key Elements of the Training and Education Plan**
The Offeror should provide a Training and Education Plan within 120 days of the subcontract award. These plans will be updated throughout the life of the project to reflect the latest content, as needed. The plans will include

- Description of available training activities that target effective use of the user environment, performance and optimization. The description should include topics, frequency, format (such as classroom training or online training, hackathons, etc.), and pricing.
- Description of collaboration with CPU and GPU vendors, other key technology providers, and NERSC staff where appropriate for the proposed training activities.

**Key Elements of the Center of Excellence for Workflow Readiness Plan**
The Offeror should provide a Center of Excellence for Workflow Readiness Plan to assist in transitioning select NERSC mission workflows to the system (e.g. NESAP focuses on workflows that contain simulations, data, and learning components) within 120 days of the subcontract award.

The plan will be updated quarterly throughout the life of the project as needed. The plan will include

- Named support staff provided by the Offeror and the CPU and GPU vendor.
- Staff training and deep-dive interactions with a set of teams.
- How application and workflow developers can begin porting and optimization activities using proposed early access systems described in Section 2
- Mutually agreed upon duration and level of effort. For example: at least 1 FTE equivalent support should be provided from the date of subcontract execution through two (2) years after final acceptance of the system.

# APPENDIX C: Facility and Site Integration Specifications

The Subcontractor will provide documentation in the **Site Preparation Facilities Plan** unless otherwise noted.

## C.1.    NERSC Facilities Overview

The NERSC-10 system will be sited at the NERSC data center in Building 59 (Shyh Wang Hall) on the Lawrence Berkeley National Laboratory campus in Berkeley, California, hereinafter referred to as Building 59. The building has four stories, with the mechanical plant occupying the lowest level. The single computer room is located on the second floor at the south end of the building.  There are two office levels located above the data center. The computer room is split into two areas: the south computing floor and the common area.

- The south computer floor accommodates the HPC and auxiliary systems.
- The common area accommodates conventional NERSC computing systems.

The data center floor is approximately 680 feet above sea level. Building 59 has a dedicated truck loading dock with a dock leveler and a freight elevator to facilitate equipment deliveries.

## C.2.    System Physical Requirements

The NERSC-10 system will be located at the north end of the south computing floor, north of the Perlmutter system at Building 59.

The approximate area of white space available for the system footprint and aisle space is 4,784 square feet (46 feet east-west, by 104 feet north-south). The Subcontractor will provide system and auxiliary racks that fit into this available space.

The Subcontractor will coordinate the optimal placement of the system with NERSC for alignment with available power and water cooling connections.

The Subcontractor will provide PDF and AutoCAD documentation that includes a complete description of the physical packaging of the system, including dimensioned drawings of individual cabinets, the definition of a scalable unit (if applicable), cooling distribution units, auxiliary racks, the row pitch, and the floor layout of the entire system.

The Subcontractor will validate the system physical size, configuration, construction, and weight are compatible with the dimension, weight, utility distribution, equipment pathway, and delivery requirements listed herein.

## C.3    Floor System & Weight Requirements

The access floor system is a conventional 24" (609.6 mm) tile system with a top elevation of +48" above the suspended concrete floor slab below. The tile pedestals are approximately 6" high and are supported on a 2" wide x 6" deep hollow structural steel (HSS) frame that is part of the facility's unique seismically base-isolated floor system. Below-floor clearance is different from conventional raised floor systems, especially at the perimeter of floor tiles where there are structural steel framing members below in all locations. This configuration limits tile cuts and utility routing near the outside ~1" of each tile.

Existing solid access floor tiles are concrete filled steel ASM FS400 units. Tiles are CISCSA load rated for a 2,000 concentrated design load, passes a 2x minimum safety factor, 6,000 pound concentrated load, 800 PSF (note structural floor live load design is for 500 PSF maximum), 200 pound impact load, 1,500 pound 10-pass rolling load, and a 1,250 pound 10,000-pass rolling load.

The Subcontractor will provide shipped and operational weights of all equipment. Equipment weights shall not exceed the capacity of the existing access floor tiles. Also provide uniform operational loads beneath cabinets, or individual foot loads, as applicable.

The Subcontractor will provide dimensioned drawings of required tile cuts based on the proposed layout of the systems for coordination and approval by the University.

## C.4    System Power & Cooling Configuration

The NERSC-10 system design will support using power and cooling connections provided below the access floor.

Power distribution from perimeter wall subpanels to computer equipment on the floor shall use underfloor cable trays supported by the HSS structural steel base isolation frames to cabinet connections at the base of the cabinets (HPC and conventional racks).

Cooling water supply and return piping is located below the access floor, distributed to equipment by 8" diameter piping manifolds. 10" diameter piping risers penetrate into the mechanical space below the computer access floor on a regular grid of approximately 12-feet. Water-cooled racks or cooling distribution units shall accommodate water feeds from below the computer access floor.

Air cooling is supplied by air handlers in the mechanical space which supply the air into a common plenum below the entire data center access floor. Conventional air-cooled racks shall be configured for consistent air flow direction to accommodate creation of hot aisle or cold aisle containment strategies.

## C.5 Seismic Requirements

LBNL is located in a seismically active area. The NERSC-10 system will be placed on a seismic isolation floor.

The Subcontractor will ensure physical integrity of the equipment (racks and all components) shall be maintained for earthquake accelerations up to 0.49g in any direction.

The Subcontractor will ensure that racks are physically interconnected near the rack base and top to provide monolithic response to each individual row of racks. Racks shall be equipped with reinforced seismic anchorage holes in the base of the frames that allows for positive, direct anchorage to the bottom of the access floor tiles using threaded fasteners and bolts. Seismic anchorage locations shall utilize four corner anchor locations at a minimum to maximize horizontal distance between holes.

The Subcontractor will provide PDF and AutoCAD drawings showing all frame dimensions and centers of gravity of all racked equipment for use in determining seismic anchorage.

**C.6 Power Requirements**
The NERSC-10 electrical system shall be compatible with 480 VAC 3-phase power fed from delta-wye transformers. Substations consist of 480V 1200A-rated distribution buckets feeding 1200A distribution boards. The Subcontractor shall provide HPC system electrical pin and sleeve connections coordinated with the University to the facility power system as 480 VAC, 3-phase (with 4 or 5 wires) or 277 VAC single phase.

The auxiliary system components are typically fed by 208 VAC 3-phase, or 208 or 120 VAC single-phase power by power distribution units provided by NERSC. Power fed at 240/415V is also possible. The Subcontractor shall coordinate layout and power needs of auxiliary systems with the University. The Subcontractor shall provide auxiliary system pin and sleeve connections coordinated with the University.

The University will provide facility power feed conductors from room perimeter wall panels underfloor to termination locations. Feeders shall use pin and sleeve receptacle termination connectors. The Subcontractor shall coordinate the type and configuration of the receptacle system with NERSC.

The NERSC-10 system should be resilient to incoming power fluctuations at least to the level guaranteed by the ITIC power quality curve and SEMI F47-0706 (Reapproved 0812) - Specification for Semiconductor Processing Equipment Voltage Sag Immunity. The Subcontractor shall provide documentation demonstrating adherence to these documents.

The Subcontractor shall ensure multi-phase power or equipment with multiple power feeds be balanced across multiple connections and phases as equally as practical.

Rate of change control: The system should provide an HPC management system alarm, via the System Scheduler, advanced warning for forecast power draw changes >1MW over a 15 minute time period, or a similar criteria based on agreement and coordination with the University. [TR-1]

All powered equipment shall be Nationally Recognized Testing Laboratories (NRTL) certified, and bear appropriate NRTL labels. NTRL certification shall be verified prior to acceptance by the University:

- For server-level products, the NTRL marks will be verified at the time of unboxing, and should be readily visible on the component case.

- For rack, custom and hybrid systems, All NTRL marks should be easily visible when equipment is installed, under power, and is in normal operating position. For equipment where this is not practical, the Vendor shall provide documentation that each item of equipment has NTRL certification prior to installation, acceptable to the University. For example, the Subcontractor may provide a list of such equipment, item serial number, physical location (e.g., rack and U location), NTRL marks (e.g. name of

testing laboratory and standard citations), and either photos of the NTRL mark, NRTL compliance documents, or reference to vendor documentation for the specific model that provides a statement of certification (e.g., provide the URL or attach the document).

All racks and system cabinets which are hard-wired, or which cannot be disconnected safely via plugs, shall provide a means for zero-voltage verification (ZVV) or provide a documented procedure for zero-voltage verification acceptable to the University and included in the specification submission.

- Finger- and tool-safe systems. All systems shall be electrically finger and tool safe, meaning that internal electrical distribution above 48V must be guarded to at least IP2X (12.5mm or larger intrusions) in accordance with IEC 60529, Degrees of Protection Provided By Enclosures (IP Code).
- Short-circuit current rating (SCCR): Provide and coordinate the short-circuit current rating capabilities for all equipment with the University to ensure the equipment is adequately provisioned for equipment protection.

For conventional, air-cooled racks:

- Utilize power distribution units (PDUs) for intra-rack power distribution.
- Use In-rack PDU strips which are acceptable to University specifications:
  - ServerTech C3WG36RL-DQJE2MT2 Switched POPS PDU (Or current generation in same model series)
- PDUs shall use 60A cables with 208V Delta 60A IEC 60309 3P+G 9h connectors.
- Utilize dual, redundant feeds with failover capability for connection to utility and UPS power.
- Cumulative peak equipment power in a single rack shall not exceed 11,500W

### C.7 Cooling Requirements
The NERSC-10 system shall be designed to operate within the following cooling conditions.

Facility cooling water (Primary Loop) conditions:

- Cooling water supply temperature range: ASHRAE W32, 5 ℃ to 32 ℃ (41 ℉ to 86 ℉), Table 5.3 of the ASHRAE Liquid Cooling Guidelines, Fifth Edition.
- Cooling water supply flow rate: Total loop flow capability of up to 16,000 GPM.
- Cooling water supply pressure: Up to 20 PSI differential pressure at the system cabinets.

Subcontractor cooling water (Secondary Loop):

- The Subcontractor shall provide operational ranges for water temperature, flow rate, and differential pressure at the CDU or cabinet level as it applies to the proposed systems.
- The Subcontractor shall develop and document a chemical management plan describing all chemical additives used for maintaining the NERSC-10 system and

cooling water property requirements in coordination with the University. The plan will include all features described in the Key Elements of the Chemical Management Plan described in Appendix C. [TR-1]

Cooling requirements for conventional air-cooled equipment:

- Conventional air-cooled equipment shall have front-to-rear air flow when installed in the rack.
- Conventional air-cooled equipment shall be compatible with a supply air meeting ASHRAE A2 ranges, Table 2.1 of the ASHRAE Liquid Cooling Guidelines, Fifth Edition. [TR-1] Ability to operate at ASHRAE A3 levels is strongly preferred.
- Conventional air-cooled system airflow requirements shall not exceed an average of 125 CFM per 1KW of load.

## C.8 Network Cabling

Network cabling (e.g., system interconnect) in Building 59 should run above floor and be integrated into the system cabinetry. There are existing cable trays that may be able to accommodate limited network cabling. NERSC-10 system network cabling should be coordinated with the University and meet the requirements:

- Permanent power, network or other cables must connect to the rear of the unit. Temporary connections for configuration or debugging are permitted on the front.
- Power, network and other cables should be neatly organized. Necessary cable management accessories are to be provided by the Subcontractor
- All network cables, wherever installed, should be source/destination labeled at both ends (see other sections for specific requirements).
- Where necessary, under floor cables shall be plenum rated and comply with NEC 300.22 and NEC 645.5.
- All network cables and fibers over 10 meters in length and installed under the floor should also have a unique serial number and dB loss data document (or equivalent) delivered at time of installation for each cable, if a method of measurement exists for cable type.

## C. 9 System Labeling

The Subcontractor will provide labels for every rack, network switch, interconnect switch, node, and disk enclosure with a unique identifier visible from the front of the and rear of the rack when the doors are open. The labels shall be high-quality plastic so that they do not fall off, fade, disintegrate, or otherwise become unusable or unreadable during the lifetime of the system. The Subcontractor will provide documentation on labeling conventions and update the documentation when changes are made.

## C.10 System Delivery and Building Pathway

The Subcontractor shall follow the requirements for system deliveries and transport of equipment through Building 59. Delivery fees and charges shall be included in the services and scope of work.

- The Subcontractor shall provide delivery drivers and personnel information with the University as required information for security and gate pass processing. Drivers shall be U.S. citizens and possess valid U.S. identification to be admitted on the LBNL site.
- Shipments and deliveries shall be made directly to LBNL Building 59
- Deliveries shall be scheduled and coordinated for specific dates and times with the University during business hours. Deliveries shall occur between 7:00 AM and 10:00 AM. For trucks in excess of 25 feet, delivery planning shall be completed at least two weeks in advance of the shipment from the origin. Final delivery dates and times shall be established at least three business days in advance for the University to ensure that security processing is given adequate review time.
- Tractor trailer delivery trucks shall be escort piloted by LBNL personnel through the Blackberry security entry gates and onto the site. LBNL flaggers will manage traffic control. Trucks need to back into the loading dock driveway. Trucks need to perform a three-point turn to leave the site with the assistance of LBNL flaggers and a pilot vehicle. Box trucks and shorter tractor trailers follow similar protocols but may be able to forego pilot vehicle and flagger requirements.
- System equipment and components shall be unloaded at the Building 59 truck loading building.
- The loading dock is an exposed, raised concrete slab structure. There is no rain cover canopy. dock, which is a single drive aisle located off of Chu Road at the north end of the
- A roll up door provides access to a concrete mechanical level and the immediately adjacent freight elevator to the second (computing) level of the building.
- A freight elevator is located at the north end of the building immediately adjacent to the loading dock roll-up door. System equipment must use this freight elevator to reach the second (Computing) level of the facility.
- The delivery path to the data center from the freight elevator is completely on raised access floor tiles. The path runs east then south through a cold shell space into the north end of the data center. A path south leads to the HPC floor area at the south end of the building. The total travel distance is approximately 250 feet.
- The pathway shall use protections, such as aluminum plates to help protect the access floor panels and moat plates. The University has a stockpile of roughly 20 aluminum plates ¼" x 48" x 96".
- The elevator cab floor shall be protected by placing an aluminum sheet inside the cab. A second aluminum sheet shall be used at each floor and placed across the elevator door opening when loading and unloading the cab to protect the elevator thresholds and avoid the load from getting caught in the gap between cab and threshold angle.
- Stabilize rolling equipment movers by using brakes or wheel chocks in the elevator to prevent the load from shifting during travel. This will reduce the risk the elevator seismic sensor will be triggered.
- The pathway shall not cross louvered or grilled tiles. Coordinate swapping out louvered and grilled tiles with solid tiles with the University, as needed.

- The Subcontractor shall prepare a documented transport and delivery plan as described in the **Key Elements of System Delivery and Installation Plan** for University review three months prior to the first scheduled system delivery.

## C.11 Building Pathway Weight and Dimension Constraints

The Subcontractor shall ensure the system design is compatible with following dimensional constraints for delivery and building pathway logistics:

- The 44" high loading dock has a leveler that is 6'-0" wide and can accommodate a vertical adjustment range of -0" to +4". The dock leveler can support a 20,000 pound load.
- The main roll up door (1201A) into the facility is electronically powered (locally) and measures 10'-7" wide by 12'-0" high.
- The freight elevator door opening measures 6'-6" wide by 9'-0" (this opening is the smallest of all openings).
- The elevator cab plan dimensions are 8'-4" wide by 8'-0" deep.
- The elevator is a C-3 freight elevator. It has a rated capacity of 7,000 pounds. The maximum allowable load in the elevator, including weight of crating, moving equipment, and personnel in the elevator cab shall be 6,300 pounds based on recommendations from the service vendor to avoid tripping the seismic motion sensor.
- There is a single interior door (2102) from the shelled space into the main computer floor. It is a double leaf door that measures 8'-0" wide by 10'-0" tall.
- The clearance between the top of the access floor and the utility systems on the bottom flange of the structural steel truss bottom chord is approximately 9'-6".

## C.12 Packaging Materials and Handling

The Subcontractor shall follow the requirements for packaging materials and handling in Building 59:

- Packaging Recycling: Packaging shall be completely recyclable and actually recycled by the Subcontractor in a reasonable and timely manner. No packaging materials shall be left behind or unaccounted for. Short-term storage of packaging materials at the Mechanical level of Building 59 can be coordinated with the University.
- Storage of materials in the computer room areas beyond a single work day, shall be coordinated and approved with the University. Per NFPA requirements, the LBNL fire marshal must approve all short-term storage plans of materials in packaging.
- Use of metal or fire-proof storage bins and containers is strongly preferred.

## C12. Safety and Training Requirements

The Subcontractor will document all emergency shutdown capabilities and internal capabilities designed to protect the system from physical damage due to electrical system faults and mechanical overheating

The Subcontractor will coordinate and document safety information for staff required to perform work on-site at LBNL with the University. The University requires

- Subcontractor Job Hazard Analysis (SJHA) to be developed in coordination with the University. This document covers the safety steps and protections taken for all work activities that are required for the system delivery and installation.
- Based on activities to be performed, the University requires LBNL training classes, or an approved, documented equivalent, reviewed by the University.
- Electrical lock-out, tag-out (LOTO): The University maintains a rigorous electrical safety program that requires formal training for any person who will perform work on or near equipment under lock-out, tag-out configuration, including coordination with University representatives for compliance with safety procedures and processes.

Hot or energized work is prohibited at LBNL.

Energization of equipment shall be made in coordination with University Qualified Electrical Workers (QEWs).

The Energization processes include University verification of appropriate voltage and phase rotation checks. Data will be shared with the Vendor.

# DEFINITIONS AND GLOSSARY

**Baseline Configuration:** A documented set of specifications for an information system, or a configuration item within a system, that has been formally reviewed and agreed on at a given point in time, and which can be changed only through change control procedures.

**Full Scale:** All of the compute nodes in the system. This may or may not include all available compute resources on a node, depending on the use case.

**Idle Power:** The projected power consumed on the system when the system is in an **Idle State**.

**Idle State:** A state when the system is prepared to but not currently executing jobs. There may be multiple idle states.

**Job Interrupt:** Any system event that causes a job to unintentionally terminate.

**Job Mean Time to Interrupt (JMTTI):** Average time between job interrupts over a given time interval on the full scale of the system. Automatic restarts do not mitigate a job interrupt for this metric.

**Node Failure:** Nodes shall be considered to have failed if a hardware problem, including soft errors, or a defect in supplied software causes the node to be unavailable, unable to operate correctly or unable to perform at established levels. A node shall be considered to have failed if a node is administratively taken offline ("admindown") to avoid erroneous operation, or by automatic action of System management software including the node health checker.

**Rolling Upgrades/Rolling Rollbacks:** A rolling upgrade or a rollback is defined as changing the operating software or firmware of a system component in such a way that the change does not require synchronization across the entire system. Rolling upgrades and rollbacks are designed to be performed with those parts of the system that are not being worked on remaining in full operational capacity.

**Software Bill of Materials:** A formal record containing the details and supply chain relationships of various components used in building software.

**System Mean Time Between Interrupt (SMTBI):** Average time between system outages over a given time interval.

**System Availability:** ((time in period – time unavailable due to outages in period)/(time in period – time unavailable due to scheduled outages in period)) * 100

**System Initialization:** The time to bring 99% of the compute resource and 100% of any service resource to the point where a job can be successfully launched.

**System Outage:** The system should be classified as down if any of the following requirements are NOT met by the System:

- 99% of compute nodes are available. This includes full health of all the node to switch links on a given node.
- At least 98.5% of all inter-switch network links are operational at full bandwidth.

(i.e., degraded links are not counted as operational.)

- At least 85% of the bandwidth out of the system from within the system is still available.
- At least 85% of the workflow environment nodes must be available to users.
- Complete benchmark applications as defined in Section 3.
- All mounted filesystems are fully operational. In other words, all users should be able to access all of their data from mounting resources to which they have access.
- Administrators are unable to access by provided APIs or ssh to the system control plane, and from there are for any reason unable to manipulate the system as needed.
- Any aspect of system monitoring or auditing is non-functional.